

**ETH login ID:**

(Please print in capital letters)

--	--	--	--	--	--	--	--	--	--	--	--

**Full name:**

---

**263-2300: How to Write Fast Numerical Code**

ETH/CS, Spring 2011

Midterm Exam

Friday, April 15, 2011

**Instructions**

- Make sure that your exam is not missing any sheets, then write your full name and login ID on the front.
- The exam has a maximum score of 100 points.
- No books, notes, calculators, laptops, cell phones, or other electronic devices are allowed.

Problem 1 (24 = 3 each)	<input type="text"/>
Problem 2 (22 = 11 + 11)	<input type="text"/>
Problem 3 (24 = 12 + 12)	<input type="text"/>
Problem 4 (30 = 10 + 6 + 12 + 2)	<input type="text"/>
<hr/>	
<b>Total (100)</b>	<input type="text"/>

## Problem 1 (24 points)

Provide **short** answers to the following questions.

A. An operation OP has the following characteristics:

- OP Latency = 7 cycles
- OP Cycles/issue = 2

Derive a lower bound on the computation time (in cycles) for the following computation

`x[0] OP x[1] OP x[2] OP ... OP x[n-1]`

Lower bound: \_\_\_\_\_ cycles

B. Consider the following program adding two vectors of length  $n$ :

```
void vectorsum(double *x, double *y, double *z, int n)
{
    int i;
    for (i = 0; i < n; i++)
        z[i] = x[i] + y[i];
}
```

If we assume a CPU that can complete 1 addition every cycle, this computation will take at least  $n$  cycles. Assume that  $x, y, z$  are all in a cache. The bandwidth between this cache and the CPU is 1 double/cycle. Derive a lower bound on the computation time (in cycles) based on this information.

Lower bound: \_\_\_\_\_ cycles

C. What is the advantage of write back caches?

D. A system with 32-bit physical addressing mode, has a 2-way set associative cache with a capacity of 64KB and block size of 32B. Derive how the address is divided:

- How many bits are used for block offset? \_\_\_\_\_
- How many bits are used for the set index? \_\_\_\_\_
- How many bits are used for tag? \_\_\_\_\_

E. For what type of cache will there never be a conflict miss?

F. Why is loop unrolling often useful in improving the performance of numerical code?

G. Consider the following c code:

```
double * func(const double N, double *x, int size){  
    for (i=1; i<size ; i++)  
        x[i]=(x[i]/size)*N;  
}
```

The previous loop can be optimized as follows:

```
double * func(const double N, double *x, int size){  
    double tmp = N/size;  
    for (i=1; i<size ; i++)  
        x[i] = x[i]*tmp;  
}
```

gcc will not perform this optimization. What could be the reason?

H. Two algorithms  $A_1$  and  $A_2$  solve the same numerical problem and involve only floating point additions and multiplications.  $A_1$  has higher performance (measured in flop/s) than  $A_2$ . Which algorithm is faster and why?

## Problem 2 (22 points)

Consider the following Matlab function implementing a fast Fourier transform (FFT). The input and output are both vectors of length  $n$ .

```
% x and y are vectors of length n
function y = fft(x)

n = length(x);

% allocate the result y which is a column vector
y = zeros(n,1);

% base case
if n == 1
    y(1) = x(1);
    return
end

m = n/2; % ignore in the cost computation

y(1:m) = x(1:m) + x(m+1:n);
y(m+1:n) = x(1:m) - x(m+1:n);

% compute_constants(m) returns a list of m constants
% the internal cost of this function can be ignored
y(m+1:n) = y(m+1:n) .* compute_constants(m);

t1 = fft(y(1:m));
t2 = fft(y(m+1:n));

y(1:2:n) = t1;
y(2:2:n) = t2;
return
```

The goal is to compute the floating point cost of this function (use the next page):

- A. Compute the number of additions (subtractions also count as additions) required for size  $n$ .

Number of additions: \_\_\_\_\_

- B. Compute the number of multiplications for size  $n$ .

Number of multiplications: \_\_\_\_\_

**Show the derivation.**

The formula  $f_k = ca^k + \sum_{i=0}^{k-1} a^i s_{k-i}$  solving  $f_0 = c$ ,  $f_k = af_{k-1} + s_k$ ,  $k \geq 1$  may be helpful.

### Problem 3 (24 points)

Consider a direct mapped cache of size 64K with block size of 16 bytes. You will calculate the miss rate for the following code using this cache. Remember that `sizeof(int) == 4`. Assume that the cache starts empty and that the local variables `i`, `j` and computations involving them take place completely within the registers.

**Note: It helps to draw the cache. Provide some detail so we see how you did it.**

A. Now consider the following code to copy one matrix to another. Assume that the `src` matrix starts at address 0 (i.e., `src[0][0]` is mapped to the first element of the first cache block) and that the `dest` matrix immediately follows it.

```
void copy_matrix(int dest[ROWS][COLS], int src[ROWS][COLS])
{
    int i, j;
    for (i = 0; i < ROWS; i++)
        for (j = 0; j < COLS; j++)
            dest[i][j] = src[i][j];
}
```

1. What is the cache miss rate if `ROWS = 128` and `COLS = 128`?  
Miss rate = \_\_\_\_\_%
2. What is the cache miss rate if `ROWS = 128` and `COLS = 192`?  
Miss rate = \_\_\_\_\_%
3. What is the cache miss rate if `ROWS = 128` and `COLS = 256`?  
Miss rate = \_\_\_\_\_%

B. Now consider the following two implementations (one is on the next page) of a horizontal flip and copy of the matrix. Again assume that the `src` matrix starts at address 0 and that the `dest` matrix immediately follows it.

```
void copy_n_flip_matrix1(int dest[ROWS][COLS], int src[ROWS][COLS])
{
    int i, j;
    for (i = 0; i < ROWS; i++)
        for (j = 0; j < COLS; j++)
            dest[i][COLS - 1 - j] = src[i][j];
    }
}
```

1. What is the cache miss rate if `ROWS = 128` and `COLS = 128`?  
Miss rate = \_\_\_\_\_%
2. What is the cache miss rate if `ROWS = 128` and `COLS = 192`?  
Miss rate = \_\_\_\_\_%

```
void copy_n_flip_matrix2(int dest[ROWS][COLS], int src[ROWS][COLS])
{
    int i, j;
    for (j = 0; j < COLS; j++)
        for (i = 0; i < ROWS; i++)
            dest[i][COLS - 1 - j] = src[i][j];
    }
}
```

1. What is the cache miss rate if ROWS = 128 and COLS = 128?  
Miss rate = \_\_\_\_\_%
2. What is the cache miss rate if ROWS = 128 and COLS = 192?  
Miss rate = \_\_\_\_\_%

## Problem 4 (30 points)

We consider a computer with a direct-mapped cache with 256 sets. Each cache block fits 16 bytes, i.e., 2 doubles. Also, the computer has 4KB pages and one TLB with 256 entries. Consider the following program which performs some computation on one column of a 2D array.

```
typedef double data[512][512]

// K is the (positive) number of iterations
// *res will contain the result; ignore in the analysis
void important_algorithm(data a, int col, int K, double *res)
{
    int i, j;
    for (i = 0; i < K; i++)
        for (j = 0; j < 512; j++)
            access_compute(a[j][col], res);
}
```

A. Determine

- (a) the number of cache misses,
- (b) the types of occurring cache misses,
- (c) the number of TLB misses

of the above program. Only consider the array `a` in this analysis. For simplicity you can assume that the column of `a` that you consider is cache-aligned, i.e., `a[0][col]` is mapped to the first element of the first cache block. Assume an empty cache and empty TLB. Again, it helps to draw the cache. Also, show some detail so we know how you did it.

B. Because the programmer was not happy with the performance, she came up with an idea to change the way the data is stored. Specifically, she introduced a new data type

```
typedef betterdata[512][513]
```

and copied using the function

```
void copy(data a, betterdata b)
{
    int i, j;
    for (i = 0; i < 512; i++)
        for (j = 0; j < 512; j++)
            b[i][j] = a[i][j];
    b[i][512] = 0;
}
```

This technique is called zero-padding. Determine the number of cache misses when the program `important_algorithm` is applied to this new data type (`betterdata` instead of `data`). Again, it helps to draw the cache. Also, show some detail so we know how you did it. `a[0][col]` is again mapped to the first element in the first cache block.

C. Since the performance of the zero-padded version was still not as good as expected, the programmer decides to buy a computer with twice the cache size (1024 doubles instead of 512). Three cache options are available. Determine in each case the number of cache misses of `important_algorithm` (still working with `betterdata`):

- (a) Twice the cache block size
- (b) Twice the number of sets
- (c) Two-way set associative

Again, it helps to draw the cache. Also, show some detail so we know how you did it.

D. With the (proper choice of) new cache the number of cache misses got indeed reduced. But by how much did the number of TLB misses get reduced and why?